

A5

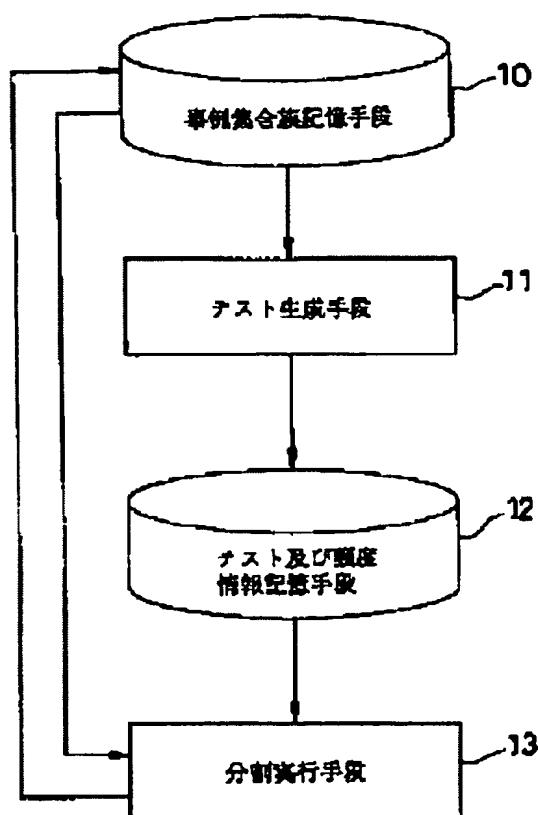
DETERMINATION TREE GENERATING METHOD

Patent number: JP9305404
Publication date: 1997-11-28
Inventor: AKIBA YASUHIRO
Applicant: NIPPON TELEGR & TELEPH CORP <NTT>
Classification:
- **international:** G06F9/44; G06F9/44
- **european:**
Application number: JP19960117890 19960513
Priority number(s):

Abstract of JP9305404

PROBLEM TO BE SOLVED: To simply express a determination tree by selecting a question, that has the most optimum training cases to be covered by the question, as a question to an attribute when there are plural optimum questions to the respective attributes having tree structure.

SOLUTION: When there are plural optimum questions to the respective attributes having the tree structure, the question having the most training case to be covered by that question is generated as the question to the attribute by a test generating means 11. That generated test is stored in a test and frequency information storage means 12, that stored test is supplied to a divided executing means 13, that divided executing means 13 selects any optimum test out of the test supplied from the test and frequency information storage means 12 and a test in the manner of ID 3, divides the test into mutually primary case sets and stores them in a case set storage means 10. By such a method, the determination tree can be simply expressed.



Data supplied from the esp@cenet database - Worldwide

THIS PAGE BLANK (USPTO)

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-305404

(43) 公開日 平成9年(1997)11月28日

(51) Int.Cl. ⁶	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 9/44	5 5 0		G 0 6 F 9/44	5 5 0 N
	5 8 0			5 8 0 E

審査請求 未請求 請求項の数 2 O L (全 6 頁)

(21) 出願番号 特願平8-117890

(22) 出願日 平成8年(1996)5月13日

(71) 出願人 000004226

日本電信電話株式会社

東京都新宿区西新宿三丁目19番2号

(72) 発明者 秋葉 泰弘

東京都新宿区西新宿三丁目19番2号 日本

電信電話株式会社内

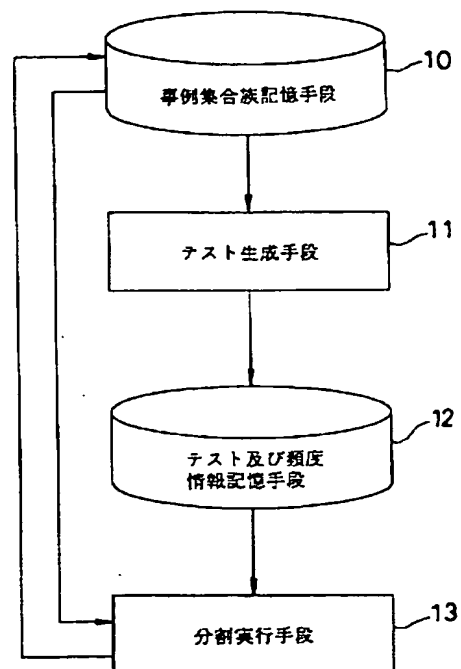
(74) 代理人 弁理士 三好 秀和 (外1名)

(54) 【発明の名称】 決定木生成法

(57) 【要約】

【課題】 生成される決定木が木構造を有する属性で表現される場合には、該木構造上でのより上位のノードに対する質問が選択される決定木生成法を提供することにある。

【解決手段】 訓練事例の表現に利用する属性の少なくとも1つが木構造を有する時、木構造を有する各属性に対して最適な質問が複数ある場合、該質問によってカバーされる訓練事例が一番多く質問を該属性に対する質問として選択し、これらの事例を説明する所望の決定木を生成する。



1

【特許請求の範囲】

【請求項 1】 いくつかの属性値と該属性値から決定されるクラスで表現された判断事例集合からルールを生成する決定木生成法であって、訓練事例の表現に利用する属性の少なくとも 1 つが木構造を有する時、木構造を有する各属性に対する最適な質問が複数ある場合、その質問によってカバーされる訓練事例が一番多い質問を該属性に対する質問として選択することを特徴とする決定木生成法。

【請求項 2】 いくつかの属性値と該属性値から決定されるクラスで表現された判断事例集合からルールを生成する決定木生成法であって、訓練事例の表現に利用する属性の少なくとも 1 つが木構造を有する時、木構造を有する各属性に対する最適な質問が複数ある場合、各質問の評価値を計算し、この計算した各質問の評価値のうち最も大きな評価値を有する質問を該属性に対する質問として選択することを特徴とする決定木生成法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、属性値とクラスで表現される判断事例集合からルールを生成する決定木生成法に関し、更に詳しくは、専門家の経験をルールとして計算機の中に取り込んで、専門家と同等の判断を可能とすることを特徴とするエキスパートシステムの構築に必要なルールを専門家の判断事例から自動的に生成する決定木生成法に関する。

【0002】

【従来の技術】 近年、専門家の経験をルールとして、計算機の中に取り込んで、専門家と同等の判断をすることを特徴とするエキスパートシステムの開発・研究が活発である。しかし、エキスパートシステム構築の大きな問題点に、専門家の知識をルールとして整理することの困難さが上げられている。この問題に対処するため、専門家の判断事例から自動的にルールを生成する方法の研究が盛んである。

【0003】 専門家の判断事例から自動的にルールを生成する代表的な方法の 1 つに、オーストラリアの Quinlan により開発された ID3 がある（「知識獲得入門・帰納学習と応用」、第 5 章帰納的推論機構、共立出版株式会社発行参照）。ID3 による事例からのルール獲得では、事例は、ある決まった属性ベクトルにより記述される。属性ベクトルの各成分の取り得る値は、属性の取り得る値である。その例を示せば、「色」とか「ボーレート」などの属性の取り得る値、{赤、緑、青}、あるいは {100, 300, 1200, 4800} が各成分の取り得る値である。ID3 では、ルールは、決定木として獲得される。

【0004】 事例集合 E に対して、事例集合 E の要素が同一のクラスに属している時の決定木は、そのクラスの名前を持つ 1 つの葉そのものである。また、事例集合 E

2

を 2 つ以上のクラスにクラス分けできるときには、適当な属性を選択することにより集合 E をそれぞれ互いに素な部分集合 E1, E2, E3, ..., En に分割することができる。ここで、Ei は E の要素の中で選択された属性の i 番目の属性値を成分に持つ要素を含む部分集合である。更に、これらの部分集合のそれぞれは、いま述べたような生成規則手順を繰り返すことによって更に処理することができる。最終的に結果は木構造となり、それぞれの葉部分にはクラスの名前が付けられている。それまでの各ノードにおいては、テストされるべき属性を規定しており、そこからの枝別れは、その属性が得られる値に対応している。

【0005】 上述した過程を説明するために、次の集合を考える。対象は、3 つの属性「背丈」、「髪の色」、「目の色」により、2 つのクラス「+」、「-」に対応づけられている。

E = 低い、ブロード、青；+
低い、黒、青；-
高い、黒、茶色；-
高い、ブロード、茶；-
高い、黒、青；-
低い、ブロード、茶；-
高い、赤、青；+
高い、ブロード、青；+

【0006】 すべての事例が同一のクラスに属しているわけではないので、テストが必要である。テストは、3 つの属性「背丈」、「髪の色」、「目の色」に対して可能である。

【0007】 例えば、髪の色に対してテストを行うと、判断木は、図 3 に示すようになる。髪の色が「黒」と「赤」の集合については、単一のクラスしか含まれていないので、これ以上分類を行う必要はない。

【0008】 一方、「ブロード」のクラスについては、複数のクラスが存在するため、同様の手順を繰り返す必要がある。

【0009】 上記の分類では、「髪の色」を分類に用いた。しかし、「背丈」、「目の色」に対しても、分類は可能であり、3 つの属性の中でどれを用いるかが実際には問題となる。ID3 では、この問題に対処するため、3 つの属性の各々を用いて分類を行った場合の分類結果のエントロピーを計算し、最もエントロピーの減少の大きな属性を用いて分類を行うことを提案している。この ID3 は、最小の分類回数で分類を終了するものではないが、ある事例を分類してクラスを決定する際に、その判断回数が（準）最小となる実用的な方法として、多用されている。上記事例集合 E に対する最終的な決定木は、図 4 に示すようなものである。

【0010】 事例を表現する属性の中には、そこに入りえる属性値が、意味的な上下関係をもつため、木構造を有するものがある。例えば、事例を表現する属性の中

3

に、属性“形”があり、属性値が、{三角形、六角形、正方形、真の楕円、円、凹}であったとすると、属性“形”は、図5に示すような木構造をもつ。このような事例を表現する属性の中に木構造を有する場合に、自動的に決定木を生成する代表的な方法の1つにAlmuallimらによる方法がある(「On handling Tree-Structured Attributes in Decision Tree learning」 in Proceedings of the 12th International Conference, Morgan Kaufmann 発行参照)。

【0011】このAlmuallimらによる方法は、木構造を有する属性について、対応する木構造Tに訓練事例を流し、訓練事例が通った木構造Tの任意のノードNに対して、“訓練事例の該属性の属性値が、木構造T上でノードNの下位ノードであるか”という質問を生成する。例えば、対象が、4つの属性「背丈」、「髪の色」、「目の色」、「好きな形」により、2つのクラスに対応づけられていて、「好きな形」が図5の木構造Tをもっている次の事例集合：

S1 = 低い、ブロンド、青、三角形；＋
 低い、黒、青、六角形；－
 高い、黒、茶、凹；－
 高い、ブロンド、茶、正方形；－
 高い、黒、青、凹；－
 低い、ブロンド、茶、凹；－
 高い、赤、青、三角形；＋
 高い、ブロンド、青、正方形；＋

を考えると、木構造をもつ属性は「好きな形」だけで、この属性に対して以下の6つのテスト：

- (1) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「凸」の下位ノードであるか？
- (2) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「正多角形」の下位ノードであるか？
- (3) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「三角形」の下位ノードであるか？
- (4) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「六角形」の下位ノードであるか？
- (5) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「正方形」の下位ノードであるか？
- (6) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「凹」の下位ノードであるか？

を生成する。最終的に、それらの質問の中で、最もエントロピーの減少の大きな質問を該属性に対する質問候補としている。

【0012】

【発明が解決しようとする課題】上記の決定木を生成する従来の技術には、大きな欠点がある。従来法により生成された決定木が木構造Tを有する属性で表現される場合には、該木構造Tの上でより上位のノードに対する質問が選択された方が、決定木が簡潔に表現され、同業者の間で広く知られているオッカムの剃刀の原理では、好ま

4

しい。しかし、従来法では、必ずしも上位のノードに対する質問が選択されるとは限らない。

【0013】本発明は、上記に鑑みてなされたものであり、その目的とするところは、生成される決定木が木構造Tを有する属性で表現される場合には、該木構造Tの上でより上位のノードに対する質問が選択される決定木生成方法を提供することにある。

【0014】

【課題を解決するための手段】上記目的を達成するため、請求項1記載の本発明は、いくつかの属性値と該属性値から決定されるクラスで表現された判断事例集合からルールを生成する決定木生成法であって、訓練事例の表現に利用する属性の少なくとも1つが木構造を有する時、木構造を有する各属性に対する最適な質問が複数ある場合、その質問によってカバーされる訓練事例が一番多い質問を該属性に対する質問として選択することを要旨とする。

【0015】請求項1記載の本発明にあつては、訓練事例の表現に利用する属性の少なくとも1つが木構造を有する時、木構造を有する各属性に対して最適な質問が複数ある場合、該質問によってカバーされる訓練事例が一番多く質問を該属性に対する質問として選択し、これらの事例を説明する所望の決定木を生成する。

【0016】また、請求項2記載の本発明は、いくつかの属性値と該属性値から決定されるクラスで表現された判断事例集合からルールを生成する決定木生成法であつて、訓練事例の表現に利用する属性の少なくとも1つが木構造を有する時、木構造を有する各属性に対する最適な質問が複数ある場合、各質問の評価値を計算し、この計算した各質問の評価値のうち最も大きな評価値を有する質問を該属性に対する質問として選択することを要旨とする。

【0017】請求項2記載の本発明にあつては、訓練事例の表現に利用する属性の少なくとも1つが木構造を有する時、木構造を有する各属性に対して最適な質問が複数ある場合、各質問の評価値を計算し、この計算した各質問の評価値のうち最も大きな評価値を有する質問を該属性に対する質問として選択し、これらの事例を説明する所望の決定木を生成する。

【0018】

【発明の実施の形態】以下、図面を用いて本発明の実施の形態について説明する。

【0019】図1は、本発明の一実施形態に係る決定木生成法を実施する装置の構成を示すブロック図である。同図に示す装置は、生成される決定木が木構造Tを有する属性で表現される場合には、該木構造Tの上、より上位のノードに対する質問が選択された決定木を生成するものであり、処理の途中の分割された事例集合の集合である事例集合族を記憶している事例集合族記憶手段10を有する。

5

【0020】この事例集合族記憶手段10に記憶されている事例集合族の各々の集合には、各事例を表現する複数の属性値と該事例が属するクラスとが記憶されている。事例集合族記憶手段10に記憶されている事例集合族は、テスト生成手段11に供給される。テスト生成手段11は、木構造を有する各属性Aについて、最適な質問が複数あれば、その質問によってカバーされる訓練事例が一番多い質問を該属性に対する質問として生成する。テスト生成手段11で生成されたテストは、テスト及び頻度情報記憶手段12に記憶される。テスト及び頻度情報記憶手段12に記憶されたテストは、分割実行手段13に供給され、分割実行手段13はテスト及び頻度情報記憶手段12から供給されたテストおよびID3流のテストから最適なテストを選択し、互いに素な事例集合を分割し、事例集合族記憶手段10に記憶する。いま述べたように生成規則手順を繰り返すことによって、最終的に結果は木構造となり、それぞれの葉部分にはクラスの名前が付けられている。

【0021】次に、図2を参照して、テスト生成手段11の処理内容の実現方法を更に詳しく説明する。テスト生成手段11は、事例集合族記憶手段10に記憶されている事例集合族から、テストを生成するものである。木構造をもつ各属性毎に図2に示す処理を行う。

【0022】まず、処理を開始すると、木構造を有する属性Aには、N個の質問候補 U_i ($i=1, 2, \dots, N$) が可能であるとし、最適な質問番号格納領域Jとカウンターjをテスト及び頻度情報記憶手段12に準備し、また各質問候補 U_i の評価値を格納する評価値記憶領域 V_i ($i=1, 2, \dots, N$) をテスト及び頻度情報記憶手段12に準備する(ステップK11, K12)。ここで、質問候補としては、回答がYESまたはNOであるものならなんでもよく、例えば、前述のAlmuallimらによる質問群を考えればよい。各質問 U_i の評価値を計算して、評価値記憶領域 V_i ($i=1, 2, \dots, N$) に格納する(ステップK13)。ここで、評価値は、決定木生成に採用される任意の評価関数により定められるもので、例えば、前述のエントロピーにより計算される。

【0023】最適な質問番号格納領域Jとカウンターjを1で初期化する(ステップK14)。カウンターjがN-1を超えないか否かをチェックし(ステップK15)、超えない場合に限り、カウンターjをカウントアップする(ステップK16)。VjがVJを超える場合か、またはVjがVJに等しくかつUjでカバーされる事例数がUjでカバーされる事例数を超える場合に限り、最適な質問番号格納領域Jをjで更新する(ステップK17, K18)。カウンターjがN-1を超えるまで、ステップK16, K17, K18の処理を繰り返す(ステップK15)。このようにして選ばれた番号Jに対応する質問を、木構造を有する属性Aに対する質問と

6

する。

【0024】図2を参照し、テスト生成手段11の具体的処理内容を説明する。対象が、4つの属性「背丈」、「髪の色」、「目の色」、「好きな形」により、2つのクラスに対応づけられていて、「好きな形」が図5の木構造Tをもっている次の事例集合を考える。

S1 = 低い、ブロード、青、三角形；+
 低い、黒、青、六角形；-
 高い、黒、茶、凹；-
 高い、ブロード、茶、正方形；-
 高い、黒、青、凹；-
 低い、ブロード、茶、凹；-
 高い、赤、青、三角形；+
 高い、ブロード、青、正方形；+

【0025】この例の場合、木構造をもつ属性は、「好きな形」だけなので、テスト生成手段の処理対象は、「好きな形」だけである。質問候補としてAlmuallimらによる質問群を採用すれば、ターゲット属性「好きな形」に対する質問候補は、以下の6つのテスト：

- (1) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「凸」の下位ノードであるか？
- (2) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「正多角形」の下位ノードであるか？
- (3) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「三角形」の下位ノードであるか？
- (4) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「六角形」の下位ノードであるか？
- (5) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「正方形」の下位ノードであるか？
- (6) 訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「凹」の下位ノードであるか？

になる。従って、この場合、質問候補数Nは6で、以下、各質問を順に U_i ($i=1, 2, \dots, 6$) とする(ステップK11)。最適な質問番号格納領域Jとカウンターjを準備し、各質問候補 U_i の評価値を格納する評価値記憶領域 V_i ($i=1, 2, \dots, 6$) を準備する(ステップK12)。ここで、評価値が、例えば、前述のエントロピーにより計算されるとすれば、各質問 U_i の評価値を計算すると、 $V_1=0.0615$, $V_2=0.0615$, $V_3=0.2442$, $V_4=0.0169$, $V_5=0.0369$, $V_6=0.0615$ となるので、これらを順に評価値記憶領域 V_i ($i=1, 2, \dots, 6$) に格納する(ステップK13)。

【0026】最適な質問番号格納領域Jとカウンターjを1で初期化する(ステップK14)。まず、カウンターjが5を超えないか否かをチェックすると(ステップK15)、超えていないのでカウンターjをカウントアップする(ステップK16)。結果として、jは2となる。V2とV1を比べると、共に0.0615で等しく、U1でカバーされる事例は、

7

低い、ブロンド、青、三角形；＋
 低い、黒、青、六角形；－
 高い、ブロンド、茶、正方形；－
 高い、赤、青、三角形；＋
 高い、ブロンド、青、正方形；＋
 の5つであり、U2 でカバーされる事例も、
 低い、ブロンド、青、三角形；＋
 低い、黒、青、六角形；－
 高い、ブロンド、茶、正方形；－
 高い、赤、青、三角形；＋
 高い、ブロンド、青、正方形；＋
 の5つであるから、Jは変化せず、1のままである（ステップK17）。jの値は2で、5を超えないので、引き続き、ステップK16、K17、K18を処理する。カウンターjをカウントアップする（ステップK16）。結果として、jは3となる。V3とV1を比べると、V3が大きい（ステップK17）。従って、Jは3で、更新される（ステップK18）。jの値は3で、5を超えないので、引き続き、ステップK16、K17、K18を処理する。カウンターjをカウントアップする（ステップK16）。結果として、jは4となる。V4とV3を比べると、V3が大きい（ステップK17）。
 【0027】従って、Jを更新せずに、ステップK15に進む。jの値は4で、5を超えないので、引き続き、ステップK16、K17、K18を処理する。カウンターjをカウントアップする（ステップK16）。結果として、jは5となる。V5とV3を比べると、V3が大きい（ステップK17）。従って、Jを更新せずに、ステップK15に進む。jの値は5で、5を超えないので、引き続き、ステップK16、K17、K18を処理する。カウンターjをカウントアップする（ステップK16）。結果として、jは6となる。V6とV3を比べると、V3が大きい（ステップK17）。従って、Jを*

8

*更新せずに、ステップK15に進む。jの値は6で、5を超えるので、図2に示す一連の処理が終了する。このようにして、選ばれた番号Jは、3なので、質問U3
 “訓練事例の属性値「好きな形」の属性値は、木構造T上でノード「三角形」の下位ノードであるか”が木構造を有する属性「好きな形」に対する質問として選択される。

【0028】

【発明の効果】以上説明したように、本発明によれば、
 10 訓練事例の表現に利用する属性の少なくとも1つが木構造を有する時、木構造を有する各属性に対して最適な質問が複数あれば、その質問によってカバーされる訓練事例が一番多い質問を該属性に対する質問として選択するので、生成される決定木が木構造Tを有する属性で表現される場合には、該木構造Tの上でより上位のノードに対する質問が選択され、結果として、決定木が簡潔に表現される。

【図面の簡単な説明】

【図1】本発明の一実施形態に係る決定木生成法を実施する装置の構成を示すブロック図である。

20 【図2】図1に示す実施形態の作用を示すフローチャートである。

【図3】事例集合の木構造の一例を示す図である。

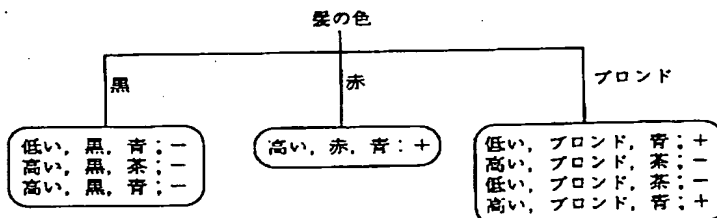
【図4】図3に示す事例集合の木構造の最終的な決定木を示す図である。

【図5】事例を表現する属性のもつ木構造の一例を示す図である。

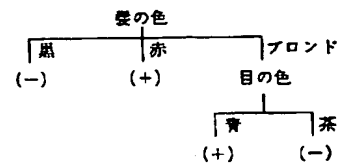
【符号の説明】

- 10 事例集合族記憶手段
- 11 テスト生成手段
- 12 テスト及び頻度情報記憶手段
- 13 分割実行手段

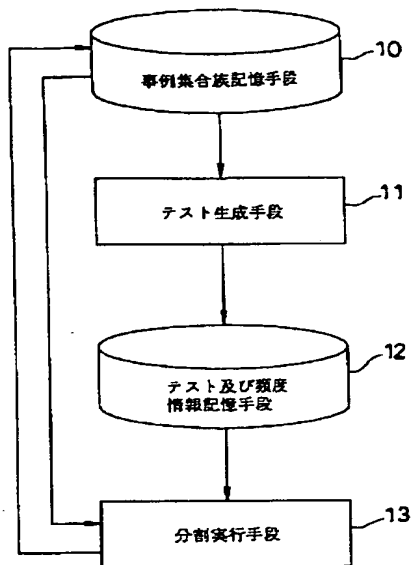
【図3】



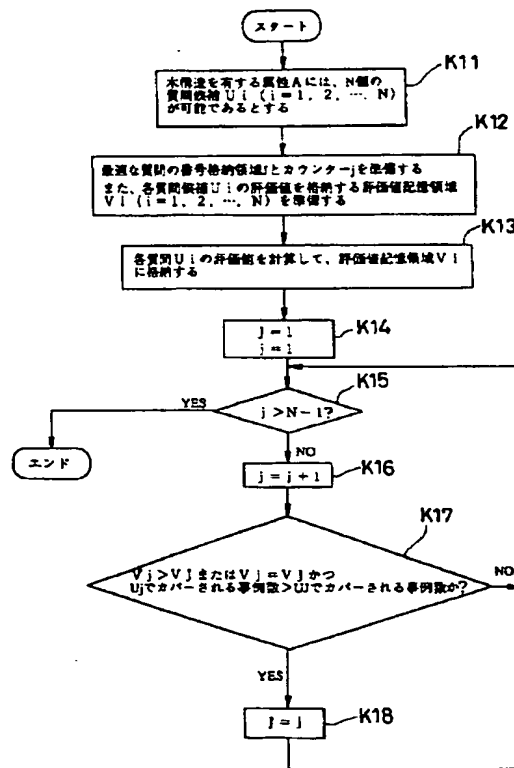
【図4】



【図1】



【図2】



【図5】

